

## Aberystwyth University

### *Concatabominations*

Siu-Ting, Karen; Pisani, Davide; Creevey, Christopher J; Wilkinson, Mark

*Published in:*  
Systematic Biology

*DOI:*  
[10.1093/sysbio/syu066](https://doi.org/10.1093/sysbio/syu066)

*Publication date:*  
2014

*Citation for published version (APA):*

Siu-Ting, K., Pisani, D., Creevey, C. J., & Wilkinson, M. (2014). Concatabominations: Identifying Unstable Taxa in Morphological Phylogenetics using a Heuristic Extension to Safe Taxonomic Reduction. *Systematic Biology*.  
<https://doi.org/10.1093/sysbio/syu066>

#### **General rights**

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

#### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400  
email: [is@aber.ac.uk](mailto:is@aber.ac.uk)

RUNNING HEAD: CONCATABOMINATIONS

TITLE: Concatabominations: Identifying Unstable Taxa in Morphological  
Phylogenetics using a Heuristic Extension to Safe Taxonomic Reduction

AUTHORS: KAREN SIU-TING<sup>1,2\*</sup>, DAVIDE PISANI<sup>2</sup>, CHRISTOPHER J. CREEVEY<sup>3</sup> AND  
MARK WILKINSON<sup>4</sup>

<sup>1</sup> *Dept. of Biology, National University of Ireland, Maynooth, Co. Kildare, Ireland*

<sup>2</sup> *School of Biological Sciences, University of Bristol, Bristol BS8 1UG, UK*

<sup>3</sup> *Institute of Biological, Environmental & Rural Sciences, Aberystwyth University,  
Aberystwyth SY23 3FG, UK*

<sup>4</sup> *Department of Life Sciences, The Natural History Museum, London SW7 5BD, UK*

*\* Corresponding author*

*Correspondence to be sent to:*

*School of Biological Sciences, University of Bristol, Bristol BS8 1UG, UK*

*E-mail: [agalychnica@gmail.com](mailto:agalychnica@gmail.com)*

KEYWORDS: Rogue taxa, consensus, resolution, taxon removal

For a variety of reasons, some phylogenetic datasets are replete with missing entries. Attitudes towards abundant missing data, specifically concerns over its potential to mislead or confound phylogenetic inferences, are varied. Thus there is a current debate on the impact of missing entries upon the accuracy of phylogenetic inferences (Wiens 2006; Lemmon et al. 2009; Philippe et al. 2011; Wiens and Morrill 2011; Roure et al. 2013). Perhaps less controversial is that individual taxa may sometimes be relatively phylogenetically unstable by virtue of limited data and extensive missing data (e.g. Wilkinson 1996; Sanderson and Shaffer 2002; Wiens 2003; Wilkinson 2003). Wilkinson (1995) developed an approach for diagnosing taxon instability due to missing data *a priori* termed safe taxonomic reduction (STR). STR allows the identification of “rogue” taxa that can be removed from a dataset safe in the knowledge that their removal will not impact upon the interrelationships that will be inferred among the remaining taxa under the parsimony criterion. The potential benefits of such deletion are reductions in numbers of optimal trees and run times and better resolved consensus summaries.

STR has been fairly widely used, mainly by palaeontologists confronted with relatively incomplete fossil taxa (see Anquetin 2012; Graf 2012; McDonald 2012; for some recent examples), but also in the context of the matrix representation with parsimony (Baum 1992; Ragan 1992) approach to supertree construction (e.g. Cardillo et al. 2004). Nonetheless STR is not always as effective as one might hope (e.g. Mannion et al. 2013). Here we present a simple heuristic method for identifying potentially unstable taxa that may be useful in cases where STR does not succeed in ameliorating all the problems caused by missing data. We illustrate the approach through application to the saurischian data of Gauthier (1986) which was previously

used to illustrate STR and thus is particularly appropriate for demonstrating the ability of the new method to achieve more than STR alone.

## THE METHOD

STR is based on the understanding that if the character states of a leaf (OTU, terminal, tip)  $w$  are a subset of those of a second leaf  $x$  (such that  $w$  and  $x$  have a pairwise-dissimilarity or p-distance of zero) then (1) there exists at least one most parsimonious tree (MPT) in which leaves  $w$  and  $x$  are a cherry (sister or adjacent taxa), and (2) removing leaf  $w$  will not alter the combinations of character states present in the data, the length of most parsimonious trees (MPTs) or relationships inferred among the remaining taxa (Wilkinson 1995). If  $w$  is similarly potentially related to multiple other leaves (e.g. to  $x$ ,  $y$ ,  $z$ , etc.) there will be multiple optimal trees that differ only in the placement of  $w$  with  $x$  or with  $y$  or with  $z$  and so on. In such cases, removing  $w$ , which adds nothing to a parsimony analysis, can be helpful in reducing numbers of equally optimal trees and improving resolution of strict consensus trees. Figure 1 gives a classification of the sorts of relations that can pertain between pairs of taxa with p-distances of zero.

Sometimes missing (*qua* limited) data seem to be a problem, as evidenced by large numbers of equally optimal trees and poorly resolved consensus trees, but STR is of limited help. In such cases there may be many pairs of leaves with p-distances of zero but, because of the distribution of missing entries, the character states of neither are a proper subset of those of the other (category D, Fig. 1). Wilkinson (1995) called such pairs of leaves "potential taxonomic equivalents that are asymmetric both ways" (we will call them D pairs) and recognised that in contrast to the other categories of taxonomic equivalence the deletion of either member of the D pair cannot be

guaranteed to be safe *a priori*. The new method we propose augments STR with a ranking of taxa intended to reflect the potential for their deletion to be safe, to substantially reduce numbers of MPTs, and to improve the resolution of strict consensus trees. Unlike STR the method is a heuristic in that the removal of candidate unstable leaves identified *a priori* by the method may not be safe, although it is not difficult to check this *a posteriori*.

The idea behind the new method is very simple. Given any D pair we can ask whether “forcing” these leaves together into a cherry on a parsimony tree would necessitate some homoplasy that is not already evident in the data. If it does not then it seems plausible that the two leaves could go together in some MPT. If one of these leaves has such a relation with many other leaves it seems plausible that this leaf will be unstable in phylogenetic analyses, which may therefore benefit from its removal.

Our approach to determining whether homoplasy is increased by forcing leaves to go together makes use of compatibility methods (e.g. Meacham and Estabrook 1985). Two characters are compatible if there is some tree on which they can both fit without any extra steps (homoplasy) and simulations have shown that compatibility decreases as homoplasy increases both for whole matrices (O’Keefe and Wagner 2001) and individual characters (Wagner 2012). We count the total number of character pairs in the data that are incompatible (Le Quesne 1969) and use this as a proxy estimate of homoplasy in the original data. We then combine the data for a D pair of leaves to make what we call a “*concatombination*” (Fig. 2), add this construct to the original data and recalculate the pairwise incompatibility. We repeat the latter for each D pair in turn. For each leaf, we define D\* as the number of times that leaf contributes to a concatombination that does not appear to increase homoplasy (i.e. does not increase the number of pairwise character incompatibilities) in the data. We

also define, for each leaf, ABC as the number of taxonomic equivalences of that leaf in the STR categories A, B or C (each of which identifies scope for *a priori* safe deletion). Taxa can be ranked based on these individual scores or their sum.

Another way of thinking about this approach is to consider that whereas no individual characters provide evidence against the hypothesis that members of a given D pair are actually the same taxon it is possible that combining their data will reveal incompatibilities (homoplasy) that provide an argument that these leaves do not belong together. Consider a data set in which all pairs of characters are incompatible. In that case adding a concatabomination can never increase the pairwise incompatibility in the matrix irrespective of whether it would entail additional homoplasy or not. In such a case  $D^*$  would be maximal for any leaves that contribute to any D pair and provides no basis for discriminating among them. Where the leaves can be ranked based on the sum of their  $D^*$  and ABC scores we envisage users safely deleting any high ranked taxa for which ABC is non-zero and then experimentally deleting the taxa with highest  $D^*$  (or  $D^* + ABC$ ) score to investigate whether this has beneficial impacts (i.e. reduction in numbers of optimal trees, increase in resolution of the strict consensus) while simultaneously checking that the deletion is safe. Removing a taxon is safe precisely when its inclusion or exclusion has no impact upon the inferred relationships of the remaining taxa, i.e., when sets of MPTs inferred with the taxon excluded or with the taxon included but subsequently pruned are identical. If tree length is insensitive to the inclusion/exclusion of a taxon this is also a good, though not infallible, indicator that it can be safely deleted (see Wilkinson 1995).

The new method has been implemented into a “*concatobominations pipeline*” in combination with STR that is available at

<http://code.google.com/p/concatabominations/>. The pipeline uses the Jeffery and Wilkinson's STR software PerlEQ v.1.0 (<http://www.molekularesystematik.uni-oldenburg.de/en/34011.html>) to find all taxonomic equivalences and Simon Harris's program COMPASS ([http://research.ncl.ac.uk/microbial\\_eukaryotes/downloads.html](http://research.ncl.ac.uk/microbial_eukaryotes/downloads.html)) to calculate incompatibility scores. The pipeline tallies the taxonomic equivalences, creates and analyses the concatabominations for every D pair and outputs D\* and ABC scores of taxa together into a file that can be loaded into Cytoscape (Shannon et al. 2003) to provide a manipulable graphical representation of the results.

#### AN EMPIRICAL EXAMPLE

We use Gauthier's (1986) morphological cladistic data for saurischians to illustrate the concatabomination approach in practice. This dataset is a much cited example of the problems of missing data in palaeontological phylogenetics (e.g., Wilkinson 1995; Kearney 2002; Norell and Wheeler 2003), having been previously used to illustrate STR (Wilkinson 1995), and comprising 17 taxa and 84 binary characters with 41% of the entries missing. Missing entries are not randomly distributed in these data but are especially concentrated in some particularly incomplete fossils taxa. Reanalysed with Paup v.4.0b10 (Swofford 2003) with branches collapsed when their maximum lengths are zero, we obtain 832,902 MPTs of 98 steps, the strict consensus of which (Fig. 3a) is disappointingly poorly resolved (with just three splits). Applied to this data set, STR identifies four taxa (*Hulsanpes*, *Liliensternus*, *Procompsognathus* and *Saurornitholestes*) that can be safely deleted *a priori*. Their deletion results in a substantial reduction in the number of MPTs (to 197, without any change in tree length) and an increase in the resolution (two additional splits) of their corresponding strict consensus tree (Fig. 3b). Note however that this

improvement of the strict consensus can be obtained through the deletion of just *Hulsanpes* and *Saurornitholestes*. Although deletions of *Liliensternus* and/or, *Procompsognathus* are both safe and reduce the number of MPTs they are not effective at increasing the resolution of the corresponding strict consensus.

Table 1 shows the data obtained from the concatabominations pipeline and Figure 4a provides a graphical representation of the same in Cytoscape with vertices representing leaves and edges connecting pairs with either (1) taxonomic equivalences in categories A, B or C (which support safe deletion rules) or (2) concatabominations that do not increase the pairwise incompatibility of the data. The two leaves with the highest  $D^*$  (*Hulsanpes* and *Saurornitholestes*) scores are also identified by traditional STR as taxa that can be safely deleted. Deletion of *Hulsanpes* alone reduces the number of MPTs for the remaining data to 45,654 without affecting tree length but does not improve (increase the number of splits in) the corresponding strict consensus. The further deletion of *Saurornitholestes* further reduces the number of MPTs to 2,758 and is sufficient to produce all the increased resolution of the consensus (from three to five splits) that can be achieved using traditional STR alone.

Beyond this the two approaches differ. Whereas STR identifies two additional taxa (*Procompsognathus* and *Liliensternus*) that can also be safely deleted, ranking based on  $D^*$  scores prompts the experimental deletion of *Coelurus*. As already noted, the deletion of *Procompsognathus* and *Liliensternus* reduces the number of MPTs (to 197) but does not further improve the strict consensus. In contrast, deletion of *Coelurus* reduces the number of MPTs to 322 and improves the resolution of the corresponding strict consensus tree by adding an additional split (Fig. 3c). Deletion of *Coelurus* does not change MPT length and the sets of trees produced from the data after its deletion are identical to the trees produced with it included but from which it



has been pruned. Thus we can be confident that the deletion of *Coelurus* is safe although it was not identified *a priori* as such by traditional STR.

We find using a graphical representation of the concatabominations pipeline output (Fig. 4), in which the degree of each vertex (leaf) represents the sum of the D\* and ABC scores, to be very useful for visualising the potential equivalence relations among the taxa and especially useful in showing how these change with the successive removal of taxa (Fig. 4b-d). Disconnected components in the graph also help identify independent sets of taxonomic equivalents (e.g., the small set including *Procompsognathus* and *Liliensternus* and the main set that contains *Hulsanpes* and *Saurornitholestes*). Rather than deleting taxa in the order suggested by the initial ranking of their scores, it makes more sense to recalculate the scores and re-rank the taxa after each deletion and this is perhaps most easily accomplished in Cytoscape. Note that after the deletion of *Coelurus* (Fig. 4d) all the taxa that were previously connected in the main set are now unconnected indicating no further potential taxonomic equivalence among those taxa.

The analysis can stop at this point because although additional safe deletions may be possible they cannot be expected to lead to sufficiently reduced numbers of MPTs such as to lead to additional splits in the corresponding strict consensus. Hence we find, *a posteriori*, that the deletions of two other taxa (*Ornitholestes* and *Microvenator*) are also safe but do not lead to any improvements of the strict consensus and are therefore quite unnecessary.

## DISCUSSION

Since its introduction, STR has been adopted, with varying degrees of success, by many phylogenetic palaeontologists as a means of identifying relatively unstable

rogue taxa that can obfuscate what analyses of the data can tell us about phylogenetic relationships of other relatively more stable taxa. It has also been applied in some supertree studies that employ matrix representations (pseudocharacter encodings) of input trees. One undoubted attraction of STR is that a taxon is deleted *a priori* only if we are certain that this deletion cannot impact upon the relationships inferred among the remaining taxa. Thus it is not like throwing away data that could have an impact on the result and is consistent with a “total evidence” philosophy.

Taxon deletion is safe whenever the sets of trees produced by (1) excluding the taxon from the data and (2) pruning it from MPTs inferred with it included are identical. In any particular case there may be useful safe taxon deletions that are not identified *a priori* using STR. Our concatabomination approach is motivated by the desire to extend or augment STR by discovering these. It is a heuristic for identifying candidate rogue taxa, the deletion of which can only be confirmed as safe *a posteriori*. It is worth noting that even the “safe” removal of taxa might impact upon branch length estimation in parametric, model-based phylogenetics and that in stratocladistics (Fischer 2008) deleting potential equivalents would be counterproductive if they are from different time intervals.

The example dataset we used to illustrate the approach served also in the development of STR and might be considered fairly well studied and understood. Thus we were surprised when application of the concatabomination approach to these data led to such a clear cut improvement over what was achievable with STR alone. The example nicely illustrates how the approach can successfully lead to additional safe taxon deletions that improve the resolution of the strict consensus tree and our understanding of what phylogenetic hypotheses are supported by the parsimonious interpretation of the data. Although the approach is heuristic, we expect that highly

ranked taxa that it identifies in practice will be the ones that most likely can be safely deleted while usefully reducing the number of MPTs.

We find the graphical representation of the results, with each taxon a vertex and edges representing potential equivalence, and the manipulation it enables to be particularly helpful. As highly connected, potentially unstable, taxa are deleted any changes in the degree of the remaining vertices and of their relative rankings will be apparent. Natural stopping points for experimental deletion are when formerly connected clusters of taxa completely separate or when connected taxa cannot be safely deleted or their safe deletion does not improve the consensus.

Recently, there has been growing interest in the detection of rogue taxa in large-scale phylogenetics mostly using purely *a posteriori* approaches (Aberer and Stamatakis 2011; Pattengale et al. 2011). Concatabominations, which sits somewhat between the pure *a priori* approach of STR and purely *a posteriori* approaches such as leaf stability (Thorley and Wilkinson 1999) or reduced consensus (Wilkinson 1994) offers another approach to this problem. That this approach can be applied to matrix representations of trees highlights its potential in diagnosing the often serious problem of ineffective overlap in broad phylogenomic (multi-gene) studies and in supertree construction (Wilkinson and Cotton 2006, Sanderson et al. 2011).

#### FUNDING

This work was supported by the Biotechnology and Biological Sciences Research Council grant BB/K007440/1 to M.W.; and by an EMBARK scholarship from the Irish Research Council awarded to K.S.

#### ACKNOWLEDGEMENTS

We thank F. Anderson, M. Charleston and P. Wagner for very helpful comments. The authors thank J. McInerney, W. Akanni and members of the Bioinformatics Lab at NUI, Maynooth for constructive discussions during the development of this method. The authors would also like to thank L. Haggerty, V. Paturyan and D. Gardner for technical help using the servers. Analyses were carried out using the computing facilities at the High-Performance Computing Centre – NUI Maynooth, Ireland and University of Bristol, UK.

#### REFERENCES

- Aberer A.J., Stamatakis A. 2011. A simple and accurate method for rogue taxon identification. *IEEE International Conference on Bioinformatics and Biomedicine*; Atlanta (GA), IEEE, p. 118-122.
- Anquetin J. 2012. Reassessment of the phylogenetic interrelationships of basal turtles (Testudinata). *J. Syst. Palaeontol.* 10:3-45.
- Baum B. 1992. Combining trees as a way of combining data sets for phylogenetic inference, and the desirability of combining gene trees. *Taxon* 41:3-10.
- Cardillo M., Bininda - Emonds R., Boakes E., Purvis A. 2004. A species - level phylogenetic supertree of marsupials. *J. Zool. (Lond.)* 264:11-31.
- Fisher D.C. 2008. Stratocladistics: integrating temporal data and character data in phylogenetic inference. *Annu. Rev. Ecol. Evol. Syst.* 39:365-385.
- Gauthier J.A. 1986. Saurischian monophyly and the origin of birds. *Mem. Calif. Acad. Sci.* 8:1-47.
- Graf J. 2012. A new Early Cretaceous coelacanth from Texas. *Hist. Biol.* 24:441-452.

- Kearney M. 2002. Fragmentary taxa, missing data, and ambiguity: mistaken assumptions and conclusions. *Syst. Biol.* 51:369-381.
- Le Quesne W.J. 1969. A method of selection of characters in numerical taxonomy. *Syst. Biol.* 18:201-205.
- Lemmon A.R., Brown J.M., Stanger-Hall K., Lemmon E.M. 2009. The effect of ambiguous data on phylogenetic estimates obtained by maximum likelihood and Bayesian inference. *Syst. Biol.* 58:130-145.
- Mannion P.D., Upchurch P., Barnes R.N., Mateus O. 2013. Osteology of the Late Jurassic Portuguese sauropod dinosaur *Lusotitan atalaiensis* (Macronaria) and the evolutionary history of basal titanosauriforms. *Zool. J. Linn. Soc.*
- McDonald A.T. 2012. Phylogeny of basal iguanodonts (Dinosauria: Ornithischia): an update. *PloS one* 7:e36745.
- Meacham C.A., Estabrook G.F. 1985. Compatibility methods in systematics. *Annu. Rev. Ecol. Syst.* 16:431-446.
- Norell M.A., Wheeler W.C. 2003. Missing Entry Replacement Data Analysis: A Replacement Approach to Dealing with Missing Data in Paleontological and Total Evidence Data Sets. *J. Vert. Paleontol.* 23:275-283.
- O'Keefe F.R., Wagner P.J. 2001. Inferring and testing hypotheses of correlated character evolution using character compatibility. *Syst. Biol.* 50:657-675.
- Pattengale N., Aberer A., Swenson K., Stamatakis A., Moret B. 2011. Uncovering Hidden Phylogenetic Consensus in Large Data Sets. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*; IEEE, p. 902-911
- Philippe H., Brinkmann H., Lavrov D.V., Littlewood D.T.J., Manuel M., Wörheide G., Baurain D. 2011. Resolving difficult phylogenetic questions: why more sequences are not enough. *PLoS Biol.* 9:e1000602.

- Ragan M. 1992. Phylogenetic inference based on matrix representation of trees. *Mol. Phylogenet. Evol.* 1:53-58.
- Roure B., Baurain D., Philippe H. 2013. Impact of missing data on phylogenies inferred from empirical phylogenomic data sets. *Mol. Biol. Evol.* 30:197-214.
- Sanderson M.J., Shaffer H.B. 2002. Troubleshooting molecular phylogenetic analyses. *Annu. Rev. Ecol. Syst.* 33:49-72.
- Sanderson M.J., McMahon M.M., Steel M. 2011. Terraces in phylogenetic tree space. *Science.* 333:448-450.
- Shannon P., Markiel A., Ozier O., Baliga N.S., Wang J.T., Ramage D., Amin N., Schwikowski B., Ideker T. 2003. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* 13:2498-2504.
- Swofford D. 2003. PAUP\*. Phylogenetic Analysis Using Parsimony (\*and Other Methods). Version 4. Sinauer Associates, Sunderland, Massachusetts.
- Thorley J.L., Wilkinson M. 1999. Testing the phylogenetic stability of early tetrapods. *J. Theor. Biol.* 200:343-344.
- Wagner P.J. 2012. Modelling rate distributions using character compatibility: implications for morphological evolution among fossil invertebrates. *Biol. Lett.* 8:143-146.
- Wiens J. 2003. Missing data, incomplete taxa, and phylogenetic accuracy. *Syst. Biol.* 52:528-538.
- Wiens J. 2006. Missing data and the design of phylogenetic analyses. *J. Biomed. Inform.* 39:34-42.
- Wiens J., Morrill M. 2011. Missing Data in Phylogenetic Analysis: Reconciling Results from Simulations and Empirical Data. *Syst. Biol.* 60:719-731.

- Wilkinson M. 1994. Common Cladistic Information and its Consensus Representation: Reduced Adams and Reduced Cladistic Consensus Trees and Profiles. *Syst. Biol.* 43:343-368.
- Wilkinson M. 1995. Coping with abundant missing entries in phylogenetic inference using parsimony. *Syst. Biol.* 44:501-514.
- Wilkinson M. 1996. Majority-rule reduced consensus trees and their use in bootstrapping. *Mol. Biol. Evol.* 13:437-444.
- Wilkinson M. 2003. Missing entries and multiple trees: Instability, relationships, and support in parsimony analysis. *J. Vert. Paleontol.* 23:311-323.
- Wilkinson M., Cotton J.A. 2006. Supertree Methods for Building the Tree of Life: Divide-and-Conquer Approaches to Large Phylogenetic Problems. In: Hodkinson T., Parnell J. editors. *Reconstructing the Tree of Life: Taxonomy and Systematics of Species Rich Taxa*. Florida: CRC Press. p. 61-75.

Table 1. Results from the concatabominations pipeline analysis of the Gauthier (1986) dataset showing numbers of D\* and ABC scores as well as the percentage of missing entries and abbreviations (Abb.) of taxon names used in the Figures.

Taxon	Abb.	% Missing			
		entries	D*	ABC	Total
<i>Hulsanpes</i>	Hul	81	7	2	9
<i>Saurornitholestes</i>	Sas	72	7	1	8
<i>Coelurus</i>	Coe	72	5	0	5
<i>Ornitholestes</i>	Ors	40	3	0	3
<i>Compsognathus</i>	Com	38	3	0	3
<i>Microvenator</i>	Mic	67	3	0	3
Ceratosauria	Cer	0	0	2	2
Deinonychosauria	Dei	6	0	2	2
Caenagnathidae	Cae	33	2	0	2
Elmisauridae	Elm	54	2	0	2
<i>Procompsognathus</i>	Pro	64	1	1	2
<i>Liliensternus</i>	Lil	48	1	1	2
Ornithomimidae	Orm	8	0	1	1
Ornithischia	Orn	0	0	0	0
Sauropodomorpha	Sau	0	0	0	0
Carnosauria	Car	2	0	0	0
Avialae	Avi	4	0	0	0



Figure 1. Hypothetical character data illustrating relations of taxonomic equivalence among pairs of taxa (after Wilkinson 1995) and the categories given in STR. Leaves  $t$  and  $u$ , which have no missing data and identical character states, are denoted actual equivalents (category A), all the other pairs have some missing data and are denoted potential equivalents. Leaves  $w$  and  $x$  have identical character data and are denoted symmetric potential equivalents (category B), all the other possible pairs (except  $t$  and  $u$ ,  $w$  and  $x$ ) are asymmetric potential equivalents. Leaves  $x$  and  $y$  are asymmetric potential equivalents both ways (category D), pairs  $y$  and  $z$ , and  $t$  and  $w$  are asymmetric all one way (categories C and E).

Figure 2. Producing a concatabomination ( $x+y$ ) for a D pair of taxa with asymmetric potential equivalence both ways. Arrows show how the concatabomination leads to a composite taxon with missing data of each original taxon replaced where possible by data from its pair. In other words, the concatabomination of a D pair is a taxon comprising the union of the character states of the D pair.

Figure 3. Strict consensus trees of MPTs for the saurischian data of Gauthier (1986) or subsets thereof showing the increase in resolution obtained by deleting taxa. a) the complete dataset (no deletions); b) after safe deletion of four taxa identified by STR; c) after deleting the highest ranked taxa identified by the Concatabominations pipeline. For abbreviations used in the trees, refer to Table 1.

Figure 4. Taxonomic equivalences inferred from the concatabominations pipeline visualised in a network with all taxa (a) and with the successive deletions of *Hulsanpes* (Hul) (b), *Saurornitholestes* (Sas) (c) and *Coelurus* (Coe) (d). Vertices represent taxa and the edges represent the type of taxonomic equivalence shared between the taxa. Vertex size is scaled to represent the amount of taxonomic equivalences a taxon has, where the bigger the vertex the more equivalences it has, hence more unstable (see scale at the bottom of figure). Types of equivalence among nodes is represented by dashed lines (types C and E) and solid lines (type D). For a complete list of abbreviations used for the taxa names refer to Table 1.

Leaf	Characters						Categories of taxonomic equivalence:
	I	II	III	IV	V	VI	
t	0	0	0	1	1	1	<div> <div>“A”</div> <div>“C” and “E”</div> </div>
u	0	0	0	1	1	1	
w	?	?	0	1	1	1	
x	?	?	0	1	1	1	<div> <div>“B”</div> <div>“D”</div> </div>
y	0	0	0	1	?	?	<div> <div>“C” and “E”</div> </div>
z	0	0	0	1	?	1	

Figure 1





Leaf	Characters					
	I	II	III	IV	V	VI
x	?	?	0	1	1	1
y	 0	 0	0	1	 ?	 ?
x + y =	0	0	0	1	1	1

Figure 2

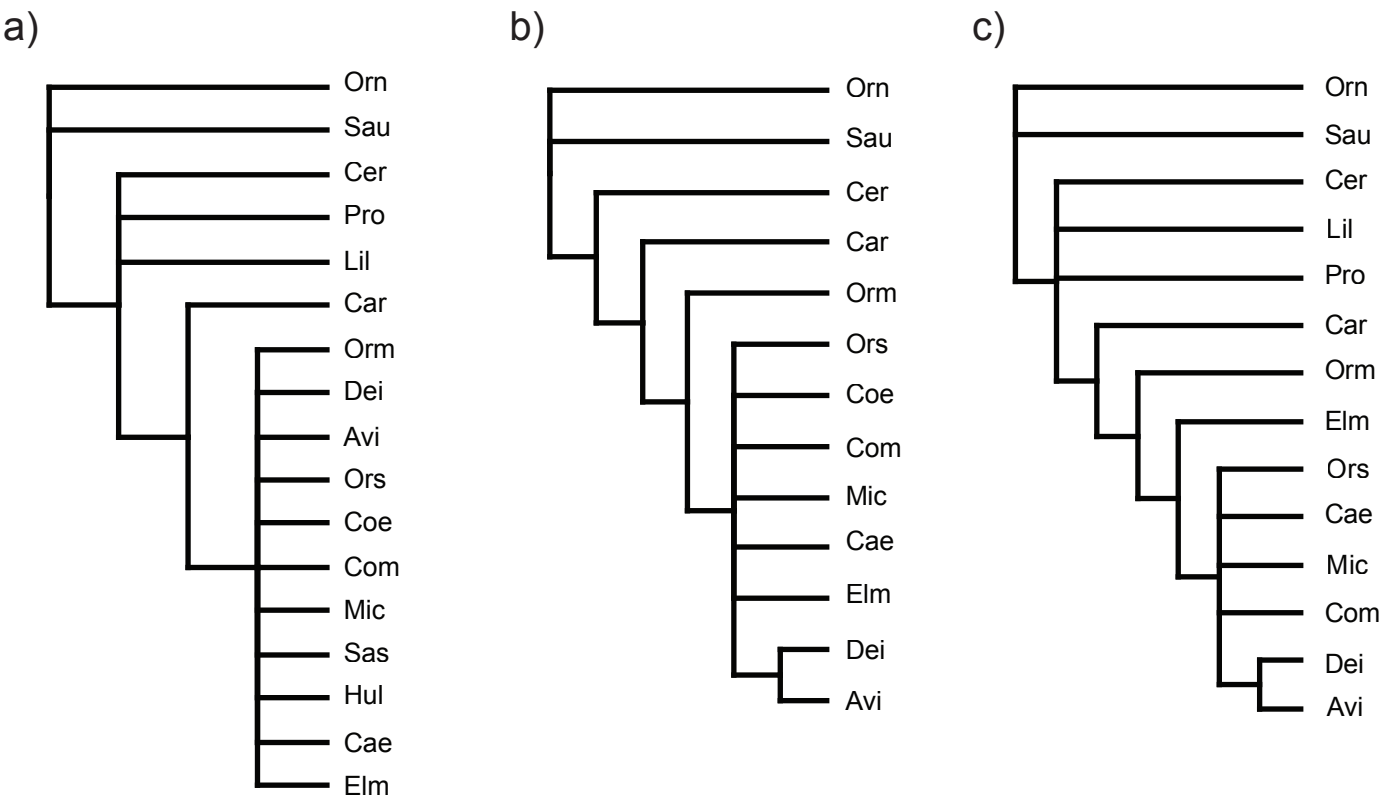


Figure 3

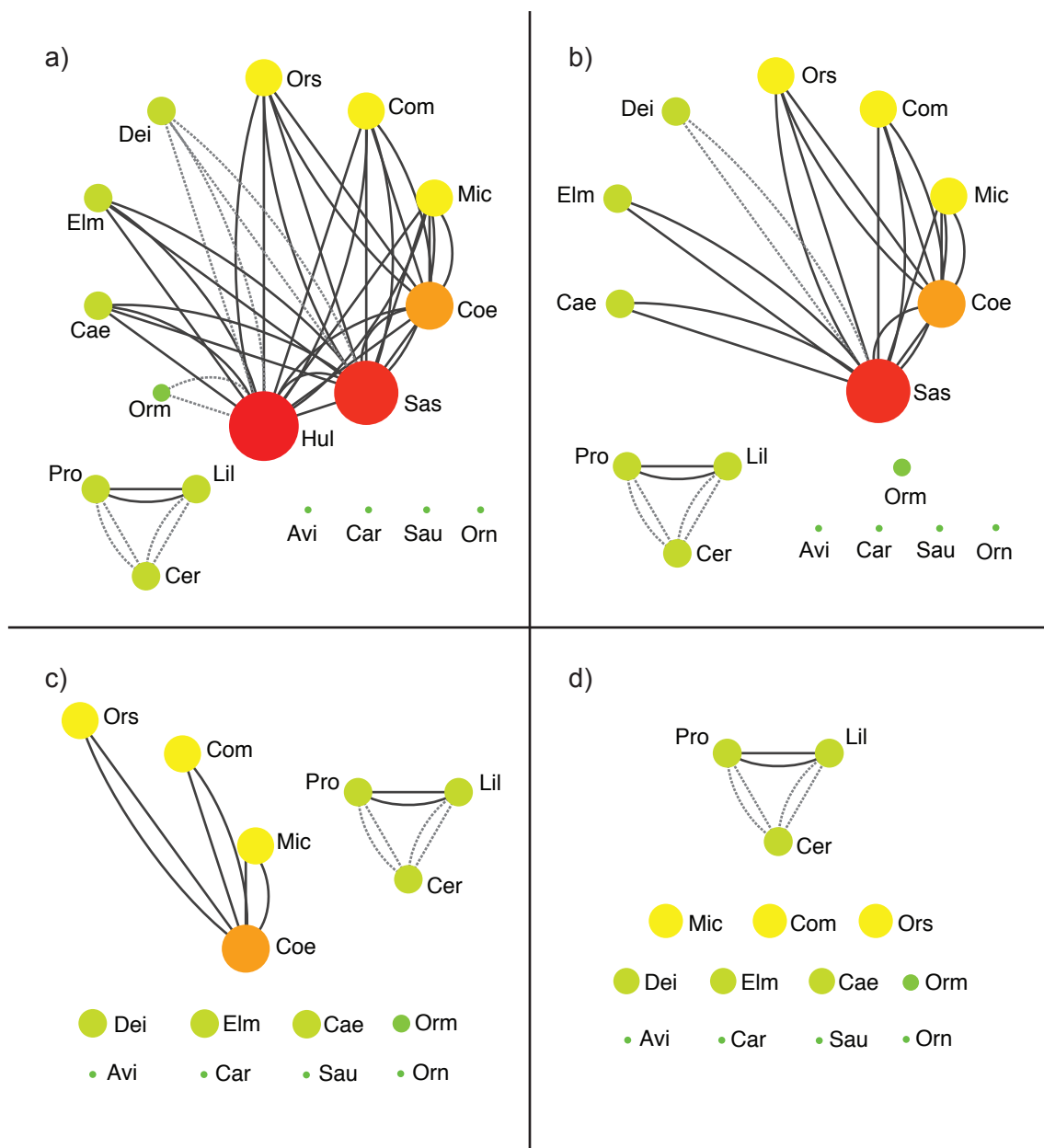


Figure 4